A Transport Protocol for Content-Based Publish/Subscribe Networks

Amir Malekpour

University of Lugano, Faculty of Informatics Advisors: Antonio Carzaniga, Fernando Pedone

Outline

- Publish-Subscribe communication paradigm
 - \triangleright Content-based
 - \triangleright Best-effort
- Problem and solution outline
- Transport Protocol
 - \triangleright FIFO ordering
 - \triangleright Reliable delivery

Publish/Subscribe Network



Publish/Subscribe Network



Publish/Subscribe Network



Events

6 / 32

RFID events 4 trillion events per day

Source: G. T. Lakshmanan et al., A stratified approach for supporting high throughput event processing applications, DEBS 2009

Events

6 / 32

- RFID events 4 trillion events per day
- Financial applications in large banks 400 million events per day

Source: G. T. Lakshmanan et al., A stratified approach for supporting high throughput event processing applications, DEBS 2009

Events

6 / 32

- RFID events 4 trillion events per day
- Financial applications in large banks 400 million events per day
- Massively multiplayer online games (MMOG)millions of events per second during peak periods.

Source: G. T. Lakshmanan et al., A stratified approach for supporting high throughput event processing applications, DEBS 2009

Publish/Subscribe

- Distriuted, asynchronous, large scale service
- A set of algorithms and potocols for routing events

Publish/Subscribe

- Distriuted, asynchronous, large scale service
- A set of algorithms and potocols for routing events

Architectures

There is a plethora of different designs and architectures

Publish/Subscribe

- Distriuted, asynchronous, large scale service
- A set of algorithms and potocols for routing events

Architectures

There is a plethora of different designs and architectures

Standards

- Java Messaging Service (JMS)
- Advanced Message Queuing Protocol (AMQP)
- Data Distribution Service (DDS)

Classifying Existing Systems

Categorizing based on many factors such as:

- Event model: topic-based, content-based
- Ordering and delivery guarantees: guaranteed, best-effort

Classifying Existing Systems

Categorizing based on many factors such as:

- Event model: topic-based, content-based
- Ordering and delivery guarantees: guaranteed, best-effort

The focus of thesis is content-based, best-effort networks

Content-based:

► An event consists of a tuple of attribute name/value

(class = "stock", symbol = "IBM", price = 125)

Content-based:

- An event consists of a tuple of attribute name/value (class = "stock", symbol = "IBM", price = 125)
- Subscriptions are a logical combination of constraints on attribute name/value: ((class = "sport") \lapha (country = "UK") \lapha (club = "Liverpool")) \lapha ((class = "stocks") \lapha (symbol = "IBM") \lapha (price < 120))

Content-based:

- An event consists of a tuple of attribute name/value (class = "stock", symbol = "IBM", price = 125)
- Subscriptions are a logical combination of constraints on attribute name/value: ((class = "sport") ^ (country = "UK") ^ (club = "Liverpool")) V ((class = "stocks") ^ (symbol = "IBM") ^ (price < 120))

Fine grained filtering

- \triangleright Numeric operators: < = >
- \triangleright String operators: substring, prefix, suffix

Content-based:

- An event consists of a tuple of attribute name/value (class = "stock", symbol = "IBM", price = 125)
- Subscriptions are a logical combination of constraints on attribute name/value: ((class = "sport") ^ (country = "UK") ^ (club = "Liverpool")) V ((class = "stocks") ^ (symbol = "IBM") ^ (price < 120))

Fine grained filtering

- \triangleright Numeric operators: < = >
- \triangleright String operators: substring, prefix, suffix

10 / 32

Guaranteed delivery and FIFO ordering

- Necessary for some applications
- Costly at run-time: higher end-to-end delay, lower throughput
- Apache ActiveMQ

10 / 32

Guaranteed delivery and FIFO ordering

- Necessary for some applications
- Costly at run-time: higher end-to-end delay, lower throughput
- Apache ActiveMQ

Best-effort systems

- Target high throughput and low end-to-end delay
 - \triangleright No "store-and-forward"
 - \triangleright No acknowledgement

10 / 32

Guaranteed delivery and FIFO ordering

- Necessary for some applications
- Costly at run-time: higher end-to-end delay, lower throughput
- Apache ActiveMQ

Best-effort systems

- Target high throughput and low end-to-end delay
 No "store-and-forward"
 - \triangleright No acknowledgement
- Less complex
 - ▷ More optimizations
 - \triangleright Hardware-assisted implementation feasible

Guaranteed delivery and FIFO ordering

- Necessary for some applications
- Costly at run-time: higher end-to-end delay, lower throughput
- Apache ActiveMQ

Best-effort systems

- Target high throughput and low end-to-end delay
 No "store-and-forward"
 - No acknowledgement
- Less complex
 - > More optimizations
 - \triangleright Hardware-assisted implementation feasible
- Datagram-like messages
 - ▷ Out-of-order deliveries
 - \triangleright Message loss
- Siena B-DRP

FIFO Ordering



FIFO Ordering





Amir Malekpour (University of Lugano) A Transport Protocol for Content-Based Publish/Subscribe Networks

11 / 32

Experiment Setup



- 8-brokers
- Each broker serves 50 clients
- ▶ 50 publishers, 350 subscribers







Time (seconds)

B-DRP throughput





B-DRP throughput

Active MQ FIFO violations

0



B-DRP throughput



Active MQ FIFO violations

0

B-DRP FIFO violations



Nearly 173,000 FIFO violations

Amir Malekpour (University of Lugano) A Transport Protocol for Content-Based Publish/Subscribe Networks

Like always...

There is a tension between desirable factors e.g., throughput vs. ordering

Like always...

There is a tension between desirable factors e.g., throughput vs. ordering

Can we mitigate this tension in best-effort systems?

Like always...

There is a tension between desirable factors e.g., throughput vs. ordering

Can we mitigate this tension in best-effort systems?

Like always...

There is a tension between desirable factors e.g., throughput vs. ordering

Can we mitigate this tension in best-effort systems?
 Not at the protocol/broker level

Like always...

There is a tension between desirable factors e.g., throughput vs. ordering

Can we mitigate this tension in best-effort systems?
 Not at the protocol/broker level
 Look at the broker network as a black box

Like always...

There is a tension between desirable factors e.g., throughput vs. ordering

- Can we mitigate this tension in best-effort systems?
 Not at the protocol/broker level
 Look at the broker network as a black box
- End-to-end "transport protocol"

Transport Protocol

▶ Not new, we've been using them since the 70s (TCP)

Transport Protocol

- ▶ Not new, we've been using them since the 70s (TCP)
- On top of a datagram (best-effort) protocol e.g.: IP
Transport Protocol

- ▶ Not new, we've been using them since the 70s (TCP)
- On top of a datagram (best-effort) protocol e.g.: IP
- Three main purposes:
 - Message ordering
 - Guaranteed delivery
 - Congestion control

Transport Protocol

- ▶ Not new, we've been using them since the 70s (TCP)
- On top of a datagram (best-effort) protocol e.g.: IP
- Three main purposes:
 - Message ordering
 - Guaranteed delivery
 - Congestion control
- But there are some fundamental differences...

Sequence Numbers



Sequence Numbers



Sequence Numbers



Amir Malekpour (University of Lugano) A Transport Protocol for Content-Based Publish/Subscribe Networks

16 / 32

FIFO Ordering

What causes FIFO violation?



What causes FIFO violation?

 $delay(m_1) > delay(m_2) + \delta$

18 / 32



What causes FIFO violation?

 $delay(m_1) > delay(m_2) + \delta$



What causes FIFO violation?

 $delay(m_1) > delay(m_2) + \delta$

$\Pr(\mathsf{FIFO violation}) = \Pr(delay(m_1) - delay(m_2) > \delta)$



What causes FIFO violation?

 $delay(m_1) > delay(m_2) + \delta$

 $\Pr(\mathsf{FIFO violation}) = \Pr(delay(m_1) - delay(m_2) > \delta)$

•
$$\underline{delay(m_1) - delay(m_2)} \sim ?$$

Observed that messages published close to each other get "swapped"

Observed that messages published close to each other get "swapped"

Observed that messages published close to each other get "swapped"

• We assume that $delay^*(m_1) = delay^*(m_2)$

19 / 32

Observed that messages published close to each other get "swapped"

- We assume that $delay^*(m_1) = delay^*(m_2)$
- ► $delay(m_1) delay(m_2) \sim (X_1 Y_1) + (X_2 Y_2) + \dots + (X_k Y_k)$ $\triangleright X_i$ and Y_i is the processing time of message at broker i

Observed that messages published close to each other get "swapped"

- We assume that $delay^*(m_1) = delay^*(m_2)$
- ► $delay(m_1) delay(m_2) \sim (X_1 Y_1) + (X_2 Y_2) + \dots + (X_k Y_k)$ $\triangleright X_i$ and Y_i is the processing time of message at broker i
- Probabilistic bounds can be found on this random variable (e.g. Bernstein inequality)
- In case of Siena B-DRP a Hypoexponential distribution fits the random variable

How does the recipient estimate the parameters?

How does the recipient estimate the parameters?

With a little help from the sender

- How does the recipient estimate the parameters?
 - With a little help from the sender
 - Sender piggybacks each message with its send time departure(m)

- How does the recipient estimate the parameters?
 - With a little help from the sender
 - Sender piggybacks each message with its send time departure(m)
- Recipient records reception time of each message arrival(m)

20 / 32

- How does the recipient estimate the parameters?
 - With a little help from the sender
 - Sender piggybacks each message with its send time departure(m)
- Recipient records reception time of each message arrival(m)
- ► x_i = delay(m_i) delay(m_j) = [arrival(m_i) - arrival(m_j)] - [departure(m_i) - departure(m_j)]

20 / 32

- How does the recipient estimate the parameters?
 - With a little help from the sender
 - Sender piggybacks each message with its send time departure(m)
- Recipient records reception time of each message arrival(m)
- ► x_i = delay(m_i) delay(m_j) = [arrival(m_i) - arrival(m_j)] - [departure(m_i) - departure(m_j)]
- Clock differences cancel out







- ▶ Is *m*¹ relevant (matching) ?
- Let's assume it is...



• We would like $Pr(FIFO \text{ violation}) \leq P_t (P_t = 0.05)$



- We would like $Pr(FIFO \text{ violation}) \leq P_t (P_t = 0.05)$
- Latch m_2 hoping for m_1 to arrive



- We would like $Pr(FIFO \text{ violation}) \leq P_t \ (P_t = 0.05)$
- Latch m₂ hoping for m₁ to arrive
- $Pr(FIFO \text{ violation}) < Pr(delay(m_1) delay(m_2) > \delta + \tau) < P_t$



- We would like $Pr(FIFO \text{ violation}) \leq P_t \ (P_t = 0.05)$
- Latch m₂ hoping for m₁ to arrive
- $Pr(FIFO \text{ violation}) < Pr(delay(m_1) delay(m_2) > \delta + \tau) < P_t$
- \blacktriangleright τ can be found using probabilistic inequalities or the Hypoexponential distribution

 Carzaniga et al. proposed an encoding scheme that encodes messages and subscriptions in Bloom filters.

22 / 32

 Carzaniga et al. proposed an encoding scheme that encodes messages and subscriptions in Bloom filters.

```
(class = "stock", symbol = "IBM", price = 301)
```

 Carzaniga et al. proposed an encoding scheme that encodes messages and subscriptions in Bloom filters.

```
(class = "stock", symbol = "IBM", price = 301)
```

 Sender piggybacks each message with encoded summary of last k messages (publication record)

m ₂	t2	0 1 1 1 0 1 1	m ₆ t ₆
m3	t3	0 1 1 0 0 0 1	class = "stock"
m4	t ₄	0 1 1 0 1 0 0	symbol = "IBM"
m5	t ₅	10100 10	price =301

 Using publication record a recipient can decide if a missing message is of interest or not (for a limited number of missing messages)

- Using publication record a recipient can decide if a missing message is of interest or not (for a limited number of missing messages)
- So if a missing message is not relevant we just don't wait for it
25 / 32

By reliability we mean going beyond best effort delivery

25 / 32

- By reliability we mean going beyond best effort delivery
- Reliable IP multicast has been studies before:
 - Reliable Multicast Protocol (RMTP) (IETF Standard)
 - Pragmatic General Multicast (PGM) (IETF Draft)
 - Scalable Reliable Multicast (SRM)
 - Cooperative message recovery

- By reliability we mean going beyond best effort delivery
- Reliable IP multicast has been studies before:
 - Reliable Multicast Protocol (RMTP) (IETF Standard)
 - Pragmatic General Multicast (PGM) (IETF Draft)
 - Scalable Reliable Multicast (SRM)
 - Cooperative message recovery
- But again...pub/sub has its own particularities

26 / 32

Two challenges:

- Two challenges:
- How to detect a missing message ?
 - ▷ Normally we would use sequence numbers (SRM)
 - Using publication record

- Two challenges:
- How to detect a missing message ?
 Normally we would use sequence numbers (SRM)
 Using publication record
- How to retrieve a missing message ?
 In SRM we would multicast a request to the multicast group
 No such thing as multicast group

- Two challenges:
- How to detect a missing message ?
 Normally we would use sequence numbers (SRM)
 Using publication record
- How to retrieve a missing message ?
 In SRM we would multicast a request to the multicast group
 No such thing as multicast group
- Using pub/sub API !

27 / 32

Volunteers subscribe to request messages
s={reliability-message=request, message=msg_id}

27 / 32

- Volunteers subscribe to request messages
 s={reliability-message=request, message=msg_id}
- A request is a published message
 r={reliability-message=request, message=msg_id}

27 / 32

- Volunteers subscribe to request messages
 s={reliability-message=request, message=msg_id}
- A request is a published message
 r={reliability-message=request, message=msg_id}
- But this is too coarse grained

- Volunteers subscribe to request messages
 s={reliability-message=request, message=msg_id}
- A request is a published message
 r={reliability-message=request, message=msg_id}
- But this is too coarse grained
- We further enhance request routing with Bloom filters
 Additional constraint/attribute that give hints about the content of messages









Amir Malekpour (University of Lugano) A Transport Protocol for Content-Based Publish/Subscribe Networks



Amir Malekpour (University of Lugano) A Transport Protocol for Content-Based Publish/Subscribe Networks

28 / 32

A pluggable module written in Java Integrates into any pub/sub system

- A pluggable module written in Java
 Integrates into any pub/sub system
- ▶ With Siena B-DRP mitigated up to 99.5% of FIFO violations

- A pluggable module written in Java
 Integrates into any pub/sub system
- ▶ With Siena B-DRP mitigated up to 99.5% of FIFO violations
- Recovered more than 70% of messages losses in very unstable networks

- A pluggable module written in Java
 Integrates into any pub/sub system
- ▶ With Siena B-DRP mitigated up to 99.5% of FIFO violations
- ▶ Recovered more than 70% of messages losses in very unstable networks
- With a small sacrifice of throughput and end-to-end delay

Future Directions

- Theoretical analysis of the protocol performance
- Better encoding schemes for publication record
 Exploiting temporal locality of events

Thank you!

Questions?